

# The lesson of causal discovery algorithms for quantum correlations: Causal explanations of Bell-inequality violations require fine-tuning

Christopher J. Wood<sup>1,2</sup> and Robert W. Spekkens<sup>3</sup>

<sup>1</sup>*Institute for Quantum Computing, Waterloo, Ontario, Canada, N2L 3G1*

<sup>2</sup>*Department of Physics and Astronomy, University of Waterloo, Waterloo, Ontario, Canada N2L 3G1*

<sup>3</sup>*Perimeter Institute for Theoretical Physics, Waterloo, Ontario, Canada N2L 2Y5*

(Dated: August 20, 2012)

An active area of research in the fields of machine learning and statistics is the development of *causal discovery algorithms*, the purpose of which is to infer the causal relations that hold among a set of variables from the correlations that these exhibit. We apply some of these algorithms to the correlations that arise for entangled quantum systems. We show that they cannot distinguish correlations that satisfy Bell inequalities from correlations that violate Bell inequalities, and consequently that they cannot do justice to the challenges of explaining certain quantum correlations causally. Nonetheless, by adapting the conceptual tools of causal inference, we can show that any attempt to provide a causal explanation of nonsignalling correlations that violate a Bell inequality must contradict a core principle of these algorithms, namely, that an observed statistical independence between variables should not be explained by fine-tuning of the causal parameters. We demonstrate the need for such fine-tuning for most of the causal mechanisms that have been proposed to underlie Bell correlations, including superluminal causal influences, superdeterminism (that is, a denial of freedom of choice of settings), and retrocausal influences which do not introduce causal cycles.

## I. INTRODUCTION

A causal relation, unlike a correlation, is an asymmetric relation that can support inferences about the consequences of interventions and about counterfactuals. The sun rising and the rooster crowing are strongly correlated, but to say that the first is the cause of the second is to say more. In particular, it says that forcing the rooster to crow early will not precipitate an early dawn, whereas causing the sun to rise early (for instance, by moving the rooster eastward), can lead to some early crowing. Nonetheless, causal structure has implications for the observed correlations and consequently one can make inferences about the causal structure based on the observed correlations. Indeed, there has been much progress in the last twenty-five years on how to make such inferences, progress that has been primarily due to philosophers and researchers in the field of machine learning and which is well summarized in the books of Pearl [1] and of Spirtes, Glymour and Scheines (SGS) [2]. Such inference schemes are known as *causal discovery algorithms*. In this article, we shall consider the question of what some prominent causal discovery algorithms have to say about the causal structure that might underlie quantum correlations, in particular those that violate Bell inequalities.

Suppose that one conducts measurements on a pair of systems that have been prepared together and then removed to distant locations such that the outcome at each wing of the experiment is outside the future light cone of the measurement choice in the other wing. Suppose further that one finds that the correlations so obtained violate Bell inequalities. If one insists on a *causal* explanation of these correlations, then it would seem that one must admit that the causes must propagate faster than the speed of light. But this is in tension with the fact that one cannot send signals faster than the speed of light. We

take this tension to be the mystery of Bell's theorem: if there are indeed superluminal causes, then why can't we use them to send superluminal signals? In this article, we will show that the principles behind causal discovery algorithms can clarify the nature of this tension. We also show that this tension persists in more exotic proposals for a causal explanations of Bell inequality violations such as superdeterminism, which is an assumption that one is not free to choose the measurement setting, and retrocausation, wherein causes propagate counter to the standard direction of time.

Our analysis will also reveal some significant inadequacies of certain existing causal discovery algorithms when applied to Bell experiments and therefore we believe that some of the expertise developed in the field of quantum foundations on causal explanations of correlations may lead to improvements in these algorithms.

The distinction between causal and inferential concepts is an instance of the distinction between ontic concepts (those pertaining to reality) and epistemic concepts (those pertaining to our knowledge of reality). Within the field of statistics, disentangling causal and inferential concepts is notoriously difficult and controversial, as is the question of when causal claims are supported by the observed correlations. In the quantum realm, where there is even less agreement about which parts of the formalism refer to ontic concepts and which refer to epistemic concepts, the problem is compounded [3]. As such, we shall try to present our analysis in a manner that does not presume any particular interpretation of quantum theory. For instance, given that different interpretations disagree on whether quantum theory implies an objective indeterminism in nature or not, we shall not presume any particular answer to this question. Instead, we simply focus on the operational predictions of the theory.

The algorithms we consider take as their input the set

of conditional independences that hold in a probability distribution over observed variables; no other feature of the probability distribution is relevant for them. In addition, we consider only algorithms that look at *statistical* independences; those that use algorithmic independences are not considered [4].

Some previous work has already considered Bell’s theorem from the perspective of causal discovery algorithms. In particular, the books by Pearl [1] and by SGS [2] comment briefly on the question. They both assert that Bell’s theorem forces a dilemma between abandoning a particular notion of locality, that there are no superluminal causal influences, and abandoning what we will call *Reichenbach’s principle*, which is the assumption that correlations are to be explained either by direct causation or a common cause. One can legitimately quibble with this conclusion on the grounds that there are other assumptions that go into Bell’s theorem: the freedom in the choice of settings and the absence of retrocausal influences, for instance. Nonetheless, we feel that this is an improvement over the standard characterization of Bell’s theorem as forcing a dilemma between abandoning locality and abandoning *realism*. It has always been rather unclear what precisely is meant by “realism”. Norsen has considered various philosophical notions of realism and concluded that none seem to have the feature that one could hope to save locality by abandoning them [5]. For instance, if realism is taken to be a commitment to the existence of an external world, then the notion of locality – that every causal influence between physical systems propagates subluminally – already *presupposes* realism.

Where previous work in causal discovery has made claims about what assumption it is best to give up in the face of Bell inequality violations, it has fallen on the side of abandoning Reichenbach’s principle.<sup>1</sup> We will take a different tack in the present work. Reichenbach’s principle will *not* be questioned. In other words, we will hold fast to the notion that *all* correlations – including those predicted by quantum theory – need to be explained causally, and we will explore what insights may be gained from causal discovery algorithms under this assumption. One reason to proceed in this manner is that the idea of explaining correlations causally appears to us to be central to the scientific enterprise. Indeed, causal hypotheses are indispensable when applying scientific theories to pragmatic ends because, unlike corre-

lations, they support inferences about the consequences of actions.

In any case, our main conclusions have not been highlighted in the previous literature on Bell’s theorem and causal discovery algorithms.

Our first conclusion is a relatively straightforward one. We note that in the case of quantum correlations for a pair of correlated systems, all correlations exhibit the following conditional independence relations among the observable variables:

1. Marginal independence of the setting variables,
2. No-signalling, that is, conditional independence of the outcome at one wing of the experiment from the setting at the opposite wing given the setting at the first wing,

and for all but a set of measure zero of experimental scenarios, these are the *only* independences. These independences characterize both the correlations that satisfy all the Bell inequalities, and the correlations that violate some Bell inequality. Therefore, if the causal discovery algorithm takes as its input not the full distribution but only the conditional independence relations that hold in the distribution (as is the case with the prominent such algorithms), then this algorithm *cannot distinguish correlations that violate Bell inequalities from correlations that satisfy Bell inequalities*. The input to such algorithms is simply too impoverished to see the difference. It follows that the causal distinctions that *do* exist between these correlations – those that are implied by Bell’s theorem – cannot be recognized by these algorithms. They will consequently make incorrect assessments of what causal structure is implied by a given set of correlations.

It is nonetheless interesting to see what the algorithms return as possible causal structures for no-signalling correlations. We look at both the case where one presumes that the settings and outcomes are the only causally relevant variables, i.e., the case of no hidden variables, and the case where one imagines that hidden variables may be causally relevant. Our main conclusion is that any causal model that can hope to explain Bell-inequality-violating correlations (or EPR correlations without recourse to hidden variables) has the feature that in order to explain the statistical independencies among the observed variables, in particular the no-signalling constraints, the model must involve a *fine-tuning* of the causal parameters, thereby violating a core principle of the best causal discovery algorithms.

So, in the end, we obtain a characterization of Bell’s theorem that is quite far from its standard characterization as a denial of “local realism”. The nebulous assumption of “realism” is replaced with Reichenbach’s principle that correlations should be explained causally. To get a contradiction, it is sufficient to supplement Reichenbach’s principle with an assumption that is rather different from Bell’s notion of local causality, namely, the assumption that the causal parameters in the model are not fine-tuned. As we shall see, the latter assumption and the

<sup>1</sup>For instance, in Ref. [6], Glymour argues that there is not even a dilemma, that one *must* abandon Reichenbach’s principle. He argues for this on the grounds that a superluminal causal influence would imply superluminal signalling, and the latter is not observed experimentally. However, this argument is incorrect because there are causal models that posit superluminal causal influences but which do not lead to superluminal signals, for instance the interpretation of quantum theory wherein the wavefunction is a complete description of reality. In such models, the Markov condition can be salvaged.

fact that there are no superluminal signals together imply the lack of superluminal causal influences, which is Bell's notion of local causality. Another advantage of this characterization of Bell's theorem is that the assumptions of Reichenbach's principle and no fine-tuning also rule out superdeterminism and certain kinds of retrocausal influences, so that these no longer consist of reasonable ways of avoiding the contradiction.

## II. CAUSAL STRUCTURES AND CAUSAL MODELS

The modern approach to the formal study of causality considers in some detail the significance of interventions and counterfactuals for defining the notion of a causal relation [1, 2]. There is a large literature on whether these sorts of definitions are adequate [7]. Although questions of this sort are relevant to a discussion of Bell's theorem, they will not be the focus of this article. We begin by describing the mathematical formalism that is relevant for describing the causal discovery algorithms in Refs. [1] and [2]. We follow the presentation of these authors.

A **causal structure** is a set of variables  $\mathbf{V}$  and a set of ordered pairs of distinct variables  $\langle X, Y \rangle$  specifying that  $X$  is a direct cause of  $Y$  relative to  $\mathbf{V}$ .

Being in a relationship of direct causation is a property that is defined relative to the set of variables being considered. If one considers a larger set which includes more variables, then what was a direct causal relation in the first set might become a mediated causal relation in the second.

Such causal structures can be represented conveniently by **directed acyclic graphs** (DAGs). A directed graph  $G$  corresponds to a set of vertices and a set of directed edges among the vertices (a vertex cannot be connected to itself). The acyclic property asserts that there are no directed paths in the graph that begin and end at the same vertex. DAGs represent causal structures in the obvious manner: every variable in  $\mathbf{V}$  is represented by a vertex, and for every pair of variables  $\langle X, Y \rangle$  where  $X$  is a direct cause of  $Y$ , there is a directed edge in the graph between the associated vertices<sup>2</sup>.

As is standard, we use the terminology of family relations in the obvious manner: if  $X$  is a cause of  $Y$ , direct or mediated, then  $X$  is said to be an *ancestor* of  $Y$ , and  $Y$  is said to be a *descendent* of  $X$ . If  $X$  is a direct cause of  $Y$ , then  $X$  is said to be a *parent* of  $Y$ . The variables in the causal structure that have no parents will be called *exogenous*, while those with parents will be called *endogenous*.

A **deterministic causal model** consists of a causal structure and a set  $\Theta$  of causal and probabilistic param-

eters. The causal parameters describe the functional relations that fix the values of every variable  $X$  given its parents  $Pa(X)$  in the causal structure, that is, for every  $X$  they describe a function  $f$  specifying  $X = f(Pa(X))$ . The probabilistic parameters specify a probability distribution over the exogenous variables, that is, a distribution  $P(X)$  for every exogenous  $X$ .

An example of a deterministic causal model is given in Fig. 1.

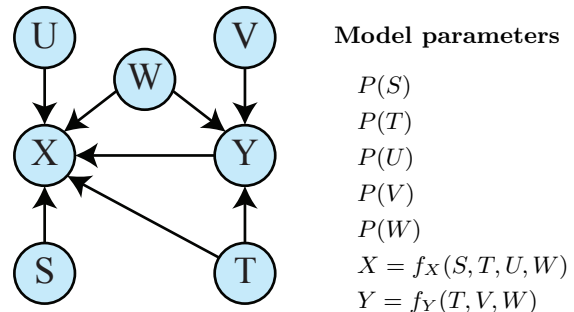


FIG. 1: An example of a deterministic causal model.

The notion of a general causal model can be explained as follows. We start with a deterministic causal model and modify it in a particular way. When an exogenous variable  $U$  is the parent of only a *single* other variable, say  $X$  (i.e. it is not a common cause of two or more variables), it is possible to eliminate  $U$  from the causal structure, and to replace the deterministic dependence of  $X$  on its original set of parents with a probabilistic dependence on its new set of parents. Specifically, if the deterministic causal model specifies that  $X = f(Pa(X))$  for some function  $f$  (here  $Pa(X)$  includes  $U$ ) then the new causal model specifies a conditional probability  $P(X|Pa'(X))$  (here  $Pa'(X)$  are the parents relative to the new causal structure, which excludes  $U$ ). Specifically, the conditional probability is defined by  $P(X|Pa(X)) = \sum_U \delta_{X, f(Pa'(X), U)} P(U)$ .

A general **causal model** consists of a causal structure and a set  $\Theta$  of causal-statistical parameters. The causal-statistical parameters specify a conditional probability distribution for every variable given its causal parents,  $P(X|Pa(X))$ . Exogenous variables have the null set for their causal parents, so that to condition on their parents is not to condition at all. Consequently, the causal-statistical parameters specify the distributions over the exogenous variables<sup>3</sup>.

An example of a general causal model is given in Fig. 2. It can be obtained from the deterministic causal model of Fig. 1 by eliminating the exogenous variables  $U$  and  $V$ .

<sup>2</sup>One can imagine more general notions of causation wherein directed cycles are allowed, but we will not consider such notions here.

<sup>3</sup>Such models are sometimes called *Markovian*. A more general sort of model, which allows bi-directed edges representing the existence of a common cause for a pair of variables, are called *semi-Markovian*.

(Note that one need not eliminate *all* exogenous variables from a deterministic causal model to obtain a nondeterministic causal model— for instance,  $S$  and  $T$  have not been eliminated in our example.)

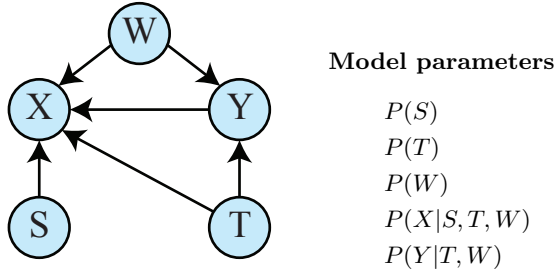


FIG. 2: An example of a causal model consisting of a causal structure, represented by a directed acyclic graph, and a set of causal-statistical parameters, specifying the probability of each variable conditioned on its parents.

Deterministic causal models are clearly a special case of causal models where all conditional probabilities correspond to deterministic functions. It is also clear that for any given causal model, one can always view it as arising from a deterministic causal model by excluding some exogenous variables. To obtain such a deterministic extension of a causal model, it suffices to add new exogenous variables as parents of every endogenous variable in the model. For the rest of the article, we will focus on the general notion of a causal model, rather than on deterministic causal models.

We pause to discuss briefly the possible interpretation of the probabilities in a causal model. One could take a Bayesian attitude towards the probabilities appearing in a causal model. In this case, the marginal probability on an exogenous variable  $U$  represents an agent's degrees of belief about  $U$ , and the conditional probability  $P(X|Pa(X))$  represents degrees of belief about  $X$  given its parents. Another possibility is to take a frequentist attitude towards the probabilities. This is arguably the position adopted by Pearl, who describes the auxiliary variables appearing in a deterministic extension of a causal model as 'unmeasurable conditions that Nature governs by some undisclosed probability function' ([1], p. 44). One could even interpret the probabilities as *propensities*, indicating an irreducible randomness in one's theory (an option that might be appealing to some when considering the possibility of explaining quantum correlations in terms of causal models). Our conclusions here will be independent of this choice <sup>4</sup>.

It is worth noting that the definition of a causal model implies that exogenous variables should be independently distributed. The idea is that we are trying to explain all correlations by a causal mechanism, so that one should

include in the model sufficiently many variables that any correlation between two variables can be deduced from the causal structure. The exogenous variables are, by definition, the variables that one takes to be uncorrelated. In this sense, the definition of a causal model incorporates Reichenbach's principle: if two variables are correlated then either one is the cause of the other, or there is a common cause. This is not an exclusive or – it could be that two variables have both a common cause and a direct causal relation between them.

Consider the following question: given a causal *model*, what sorts of correlations can be observed among the variables? Clearly, there is a set of joint distributions that are possible, depending on the causal-statistical parameters that we add to the causal structure to get a causal model.

Consider the example from Fig. 2. It is clear that the causal model predicts that the joint distribution over all the variables should be

$$P(X, Y, S, T, W) = P(W) P(S) P(T) P(Y|T, W) \times P(X|Y, S, T, W). \quad (2.1)$$

In general, a causal model with variables  $\mathbf{V} \equiv \{X_1, \dots, X_n\}$  predicts a joint distribution of the form

$$P(X_1, \dots, X_n) = \prod_{i=1, \dots, n} P(X_i | Pa(X_i)). \quad (2.2)$$

Essentially, one multiplies together the conditional probabilities for every variable given its parents, all of which are specified by the causal model. For a DAG that is not a complete graph (i.e not every pair of nodes is connected by an edge), the probability distributions that it supports are a subset of the possible distributions over those variables.

We now turn to another question: what properties do all distributions consistent with a given causal structure have in common? In other words, what are the features of the joint probability distribution that depend *only* on the causal structure and not the causal-statistical parameters? Conditional independence (CI) relations are an example of such properties, and they are the ones that most causal discovery algorithms focus upon.

Recall that variables  $X$  and  $Y$  are conditionally independent given  $Z$ , denoted

$$(X \perp Y \mid Z)$$

if any of the following three equivalent conditions hold

1.  $P(X|Y, Z) = P(X|Z) \forall P(Y, Z) > 0$ ,
2.  $P(Y|X, Z) = P(Y|Z) \forall P(X, Z) > 0$ ,
3.  $P(X, Y|Z) = P(X|Z)P(Y|Z) \forall P(Z) > 0$ .

An intuitive account of each of these conditions is as follows: In the context of already knowing  $Z$ , (1) learning  $Y$  teaches you nothing about  $X$  (i.e.  $Y$  teaches you nothing more about  $X$  than what you already could infer from

<sup>4</sup>Although we ultimately favor the Bayesian interpretation.

knowing  $Z$ ), (2) learning  $X$  teaches you nothing about  $Y$ , and (3)  $X$  and  $Y$  are uncorrelated. Note that *marginal independence* of  $X$  and  $Y$ , where  $P(X, Y) = P(X)P(Y)$ , is simply conditional independence where the conditioning set is the null set.

The definition of conditional independence implies that certain logical inferences hold among CI relations. In other words, a set of CI relations need not be logically independent. In particular, the *semi-graphoid axioms* specify some inferences that can be drawn among CI relations. They are:

$$\begin{aligned} \text{Symmetry:} & \quad (X \perp Y | Z) \Leftrightarrow (Y \perp X | Z) \\ \text{Decomposition:} & \quad (X \perp YW | Z) \Rightarrow (X \perp Y | Z) \\ \text{Weak Union:} & \quad (X \perp YW | Z) \Rightarrow (X \perp Y | ZW) \\ \text{Contraction:} & \quad (X \perp Y | Z) \text{ and } (X \perp W | ZY) \\ & \quad \Rightarrow (X \perp YW | Z) \end{aligned}$$

Any set of variables can be considered as a new variable, so each of the variables  $X, Y, W$  and  $Z$  appearing in the axioms should be understood as possibly representing a set of variables. These axioms are quite intuitive. Decomposition, for instance, states that if, in the context of knowing  $X$ , learning  $W$  and  $Y$  teaches you nothing about  $U$ , then learning  $W$  *alone* teaches you nothing about  $U$ .

Note that if one wants to specify *all* the CI relations that hold for a given probability distribution, it suffices to specify a *generating set*, defined to be a set from which the rest can be obtained by the semi-graphoid axioms. In this paper, the conditional independence relations will typically be specified by a generating set.

With these tools in hand, we can now discuss the central result concerning what properties of a joint probability distribution can be inferred from the causal structure.

**Theorem 1 (Causal Markov condition)** *In the joint distribution induced by a causal structure, every variable  $X$  is conditionally independent of its nondescendants given its parents,*

$$(X \perp Nd(X) | Pa(X)).$$

This result follows from Eq. (2.2) because

$$\begin{aligned} & P(X | Pa(X), Nd(X)) \\ &= \frac{P(X, Pa(X), Nd(X))}{P(Pa(X), Nd(X))}, \\ &= \frac{P(X | Pa(X)) \prod_{Y \in Pa(X), Nd(X)} P(Y | Pa(Y))}{\prod_{Y \in Pa(X), Nd(X)} P(Y | Pa(Y))}, \\ &= P(X | Pa(X)). \end{aligned} \quad (2.3)$$

The causal Markov condition implies a CI relation for every variable that is not exogenous in the causal structure. One can then infer additional CI relations from these by the semi-graphoid axioms.

To see these ideas in action, consider again the example from Fig. 2. It turns out that  $(Y \perp S | T)$  for this causal

structure, as we now demonstrate. Applying the causal Markov condition to  $Y$ , one infers that  $(Y \perp XS | WT)$ . Applying it to  $W, S$  and  $T$  one infers  $(W \perp ST), (S \perp WT)$  and  $(T \perp WS)$  respectively. By the decomposition axiom,  $(Y \perp XS | WT)$  implies  $(Y \perp S | WT)$ . From the contraction axiom,  $(Y \perp S | WT)$  and  $(S \perp WT)$  imply  $(S \perp YWT)$ . Finally, from weak union we obtain  $(S \perp YW | T)$  and then from decomposition again we have  $(S \perp Y | T)$ , which is equivalent by symmetry to  $(Y \perp S | T)$ .

We see that it can be rather laborious to infer CI relations from the causal Markov condition and the semi-graphoid axioms. Fortunately, there is a graphical criterion for identifying such relations, known as *d-separation* [1]. We will not dwell on this notion here, but we present a brief introduction in App. A.

Note that in addition to the CI relations that are implied by the causal structure, the particular causal-statistical parameters may imply other such relations. Such additional CI relations are problematic for causal discovery algorithms, as we shall see.

### III. CAUSAL DISCOVERY ALGORITHMS

We have described the correlations that are possible for a given causal structure. Causal discovery algorithms seek to solve the inverse problem: starting from correlations among observed variables, can one infer which causal structures might account for these correlations? Researchers in this area have indeed devised some schemes for narrowing down the set of causal structures that can yield a natural explanation of the correlations, wherein the notion of naturalness at play is one that we shall make explicit shortly. The algorithms look to the conditional independences among the variables to infer information about the causal structure.

In general, causal discovery algorithms may be applied directly to experimental data and in this case one needs to deal with the subtle issue of how to infer conditional independence relations from a finite sample of a probability distribution. However, we are here going to apply the causal discovery algorithms directly to the distributions prescribed by quantum theory, so we needn't worry about this subtlety.

It is worth reviewing a few basic facts about the output of causal discovery algorithms. First of all, two different causal structures might support precisely the same probability distributions, so that observation of one of these distributions necessarily leaves one ignorant about which causal structure is at play. As an example, for three variables, the three causal structures show in Fig. 3 all support the same set of probability distributions – those wherein  $A$  and  $B$  are conditionally independent given  $C$  (these are the DAGs wherein  $A$  and  $B$  are d-separated given  $C$ ). (The general conditions under which two causal structures are observationally equivalent is given by theorem 1.2.8 in Ref. [1].)

It follows that causal discovery algorithms will nec-

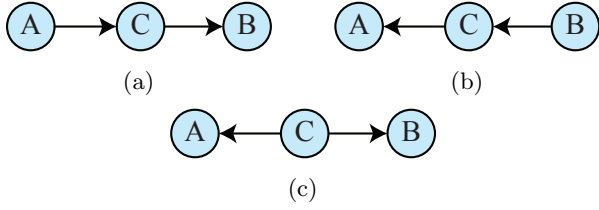


FIG. 3: The three causal models consistent with the CI relation ( $A \perp B | C$ )

essarily sometimes yield an equivalence class of causal structures. When this occurs, additional information is required if one is to narrow the causal structure down to a unique possibility, for instance information about the temporal order of some of the variables.

Despite this, one can often narrow down the field of causal possibilities significantly. To get a feeling for how this works, it is useful to start with a very simple example. Suppose that one has three binary-valued variables, denoted  $A$ ,  $B$  and  $C$ . Suppose further that the joint distribution over the triple,  $P(A, B, C)$  is such that

$$\begin{aligned} (A \perp B) \quad \text{i.e.} \quad P(A, B) &= P(A)P(B), \\ (A \not\perp C) \quad \text{i.e.} \quad P(A, C) &\neq P(A)P(C), \\ (B \not\perp C) \quad \text{i.e.} \quad P(B, C) &\neq P(B)P(C). \end{aligned} \quad (3.1)$$

What is the *natural* causal explanation for this sort of correlation? It is as shown in Fig. 4. The marginal independence of  $A$  and  $B$  is explained by their being causally independent.



FIG. 4: The *natural* causal model for the set of CI given in Eq. (3.1)

However, there are other possible causal explanations, such as the one given in Fig. 5. The reason this is a possible explanation is because there are two causal mechanisms by which  $A$  and  $B$  could become correlated, and it *could be* that the two types of correlations combine in such a way as to leave  $A$  and  $B$  marginally independent. For this to happen, however, the parameters in the causal model cannot be chosen arbitrarily and it is in this sense that the explanation is less natural than the one provided by Fig. 4.

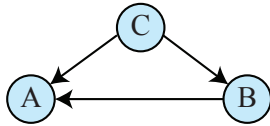


FIG. 5: An *unnatural* causal model for the set of CI given in Eq. (3.1)

An example helps to make all of this more explicit.

We adopt the following notational convention (inspired by the representation of mixtures in quantum theory)

$$\begin{aligned} P(A) &= [x] \text{ means } P(A = x) = 1, \\ P(A, B) &= [x][y] \equiv [xy] \text{ means } P(A = x, B = y) = 1. \end{aligned}$$

Consider the following joint distribution, which has the dependences described in Eq. (3.1),

$$P(A, B, C) = \frac{1}{4}[000] + \frac{1}{4}[010] + \frac{1}{4}[100] + \frac{1}{4}[111]. \quad (3.2)$$

We can easily verify that

$$P(A, B) = \left( \frac{1}{2}[0] + \frac{1}{2}[1] \right) \left( \frac{1}{2}[0] + \frac{1}{2}[1] \right),$$

so that  $A$  and  $B$  are indeed marginally independent. We also have

$$P(A, C) = P(B, C) = \frac{1}{2}[00] + \frac{1}{4}[10] + \frac{1}{4}[11],$$

so that  $A$  and  $C$  are marginally dependent, as are  $B$  and  $C$ .

The natural explanation is achieved by assuming that the causal structure is as given in Fig. 4, and the priors over the exogenous variables and the conditional probabilities for the endogenous variables are as follows:

$$\begin{aligned} P(A) &= \frac{1}{2}[0] + \frac{1}{2}[1], \\ P(B) &= \frac{1}{2}[0] + \frac{1}{2}[1], \\ P(C|A, B) &= [A \cdot B], \end{aligned}$$

where  $A \cdot B$  denotes the product of the values of  $A$  and  $B$ . Thus in this causal model,  $A$  and  $B$  are each chosen uniformly at random, and  $C$  is obtained as their product (equivalently, the logical AND of  $A$  and  $B$ ). One can easily verify that  $P(A)P(B)P(C|A, B)$  yields the distribution of Eq. (3.2).

The alternative explanation assumes the causal structure of Fig. 5, with parameters

$$\begin{aligned} P(C) &= \frac{3}{4}[0] + \frac{1}{4}[1], \\ P(B|C = 0) &= \frac{2}{3}[0] + \frac{1}{3}[1], \\ P(B|C = 1) &= [1], \\ P(A|B = 0, C = 0) &= \frac{1}{2}[0] + \frac{1}{2}[1], \\ P(A|B = 1, C) &= [C]. \end{aligned}$$

(We need not specify  $P(A|B = 0, C = 1)$  because  $P(B = 0, C = 1) = 0$ .) The joint distribution one obtains is again that of Eq. (3.2).

The difference between the two explanations becomes clear when we vary the parameters. If we change the parameters in the first model, for instance to

$$\begin{aligned} P(A) &= w[0] + (1 - w)[1], \\ P(B) &= w'[0] + (1 - w')[1], \\ P(C|A, B) &= w''[AB] + (1 - w'')[A \oplus B], \end{aligned}$$



where  $\oplus$  denotes addition modulo 2, then the joint distribution is no longer of the form of Eq. (3.2), but it is still true that  $A$  is independent of  $B$ , while  $A$  and  $C$  are dependent, and  $B$  and  $C$  are dependent. On the other hand, modifications to the parameters in the second model do not preserve the pattern of dependences and independences among  $A$ ,  $B$  and  $C$ .

The first causal structure explains the pattern of statistical dependences and independences in a manner that is robust to changes in the parameters of the causal model, whereas the second causal structure does not. Causal discovery algorithms therefore favour the first model over the second.

In the example we have used, all of the variables in the causal model were observed variables. In general (and especially in a quantum context), one might only observe a subset of the variables that are part of the causal model. Even in this case, however, one should prefer those causal models wherein the conditional independences in the probability distribution over the observed variables are stable to changes in the causal-statistical parameters.

This is the main assumption of the causal discovery algorithms, usually called *faithfulness* [2] or *stability* [1]. It is the key assumption in our analysis, so we highlight it:

**Faithfulness:** The probability distribution induced by a causal model  $M$  (over the variables in  $M$  or some subset thereof) is faithful if its conditional independences continue to hold for any variation of the causal-statistical parameters in  $M$ .

In other words, causal discovery algorithms assume that any conditional independences in the observed statistics are not a result of *fine-tuning* of the causal-statistical parameters. All independences should be a consequence of the causal structure alone. For almost any probability density over the parameter space, the parameter choices that can explain the statistical dependences in question within the *unnatural* causal structures will be found to have measure zero.

The second major assumption of causal discovery algorithms is an appeal to Occam’s razor, an assumption that one should favour the most simple or most minimal model that explains the statistics. Again, it can be applied both for the case where the observed variables are all the variables in the causal model, or the case where they are some subset thereof.

A causal model  $M$  will be said to *simulate* another causal model  $M'$  on a set of variables  $V$  if for every choice of causal-statistical parameters on  $M'$ , there is a choice of causal-statistical parameters on  $M$  such that  $M$  yields the same distribution over  $V$  as  $M'$  does. We can now define the assumption of minimality.

**Minimality:** Given two causal models  $M$  and  $M'$  that induce a given probability distribution over a set of observed variables  $V_O$  (in general a subset of the variables postulated by each causal model), if  $M'$  can simulate  $M$  on  $V_O$  but  $M$  cannot simulate  $M'$  on  $V_O$ , then  $M$  is

preferred to  $M'$  as a causal explanation of the probability distribution over  $V_O$ .

At first sight, it might seem odd to prefer  $M$  over  $M'$  given that  $M$  is consistent with fewer distributions over  $V$  than  $M'$  is. But the fact that  $M$  can explain *less* than  $M'$  implies that  $M$  is *more falsifiable* than  $M'$ , and in the version of Occam’s razor espoused by causal discovery algorithms, the degree of falsifiability is the figure of merit that one seeks to optimize. More falsifiable theories are to be preferred because, in Pearl’s words, “they provide the scientist with less opportunities to overfit the data “hind-sightedly” and therefore command greater credibility if a fit is found” ([1], p. 49). It follows that a causal model is deemed most simple if it has the *least* expressive power, while still doing justice to the observed probability distribution. Note that  $M$  might be preferred to  $M'$  as a causal explanation of the probability distribution over  $V_O$  even though  $M$  may require *more* latent variables and/or *more* causal arrows than  $M'$ ; “the preference for simplicity [...] is gauged by the expressive power of a structure, not by its syntactic description.” ([1], p. 46). We will see some examples of the consequences of the assumption of minimality shortly.

It is worth remembering that causal discovery algorithms are fallible. They are best considered a heuristic, an inference to the best explanation. Indeed, Pearl likens the faithfulness assumption in causal discovery to the following kind of inference: you see a chair before you and infer that there is a single chair rather than two chairs positioned such that the one hides the other ([1], p. 48). The task of causal discovery can be understood as “an inductive game that scientists play against Nature” ([1], p. 42).

#### A. Example of causal discovery assuming no latent variables

Variables that are not observed but which are causally relevant are called *latent* variables, or *hidden* variables. In this section, we assume that the observed variables are the *only* causally relevant variables, i.e. that there are no hidden variables. We look at a particular example of how faithfulness can help to determine candidate causal structures from a pattern of dependences in this case. The scheme is equivalent to the one introduced by Wermuth and Lauritzen [8].

Suppose one is interested in answering the question “Does smoking cause lung cancer?” For each member of a population of individuals, the value of a variable  $S$  is known, indicating whether the individual smoked or not, and the value of a variable  $C$  is known, indicating whether they developed cancer or not. Suppose a correlation between  $S$  and  $C$  is observed. Furthermore, suppose that one also has access to a third variable  $T$ , indicating whether the individual had tar in their lungs or not, and suppose that it is found that  $S$  and  $C$  are conditionally independent given  $T$ . In other words, after conditioning

on whether or not there is tar in the lungs, smoking and lung cancer are no longer correlated. Finally, imagine that these three variables are assumed to be the only causally relevant ones (we will consider the alternative to this assumption further on). What causal structure is natural given the observed conditional independence relation? Because we wish to make it very clear how these algorithms work, we will not simply specify what causal structure they return. Instead, we will look “under the hood” of these algorithms.

We begin by considering every possible hypothesis about the causal ordering. A causal ordering of variables is an ordering wherein causal influences can only propagate from one variable to another if the second is higher in the order than the first.

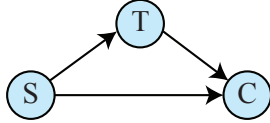


FIG. 6: The most general DAG for the causal ordering  $S < T < C$ .

For instance, consider the causal ordering  $S < T < C$ . The most general causal structure consistent with such an ordering is given in Fig. 6. To get a causal model, we need to supplement this with conditional probabilities of every variable given its parents, that is,  $P(S)$ ,  $P(T|S)$ , and  $P(C|T, S)$ . The joint distribution that this model defines is simply

$$P(S, T, C) = P(S)P(T|S)P(C|T, S).$$

Given that *any* distribution can be decomposed in this form, by choosing the conditional probabilities appropriately, we can model *any* joint distribution  $P(S, T, C)$ . But now we make use of the additional information we have about the joint distribution, namely that  $(S \perp C|T)$ . This implies that we can take the parameters in the causal model to be such that  $P(C|T, S) = P(C|T)$ , so that the joint distribution can be written as

$$P(S, T, C) = P(S)P(T|S)P(C|T),$$

and we can drop the causal arrow from  $S$  to  $C$ , so that the underlying causal structure is simply given by Fig. 7.

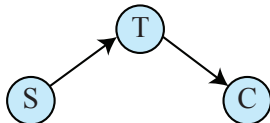


FIG. 7: DAG that captures  $(S \perp C|T)$  for the causal ordering  $S < T < C$ .

This simplified causal structure cannot generate an arbitrary probability distribution, but it *can* generate one

wherein  $(S \perp C|T)$ . It is a candidate for the true causal structure.

One then simply repeats this procedure for every possible choice of the causal ordering. For instance, for the ordering  $C < T < S$ , the most general causal structure is the one shown in Fig. 8. The decomposition of the joint

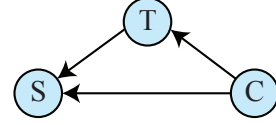


FIG. 8: The most general DAG for the causal ordering  $C < T < S$ .

probability corresponding to this causal structure is

$$P(S, T, C) = P(C)P(T|C)P(S|C, T),$$

but the constraint  $(S \perp C|T)$  implies that one substitute  $P(S|C, T) = P(S|T)$  in the causal model. Therefore, by the assumption of minimality, we drop the causal arrow from  $C$  to  $S$ , yielding a causal structure of the form given in Fig. 9. So this is another possible causal structure.

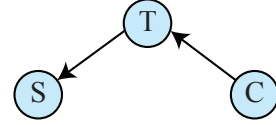


FIG. 9: DAG that captures  $(S \perp C|T)$  for the causal ordering  $C < T < S$ .

Sometimes different causal orderings lead to the same causal structure, for instance, the orderings  $T < S < C$  and  $T < C < S$  both yield the structure given in Fig. 10.

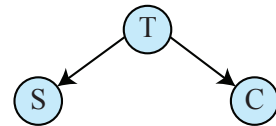


FIG. 10: DAG that captures  $(S \perp C|T)$  for the causal orderings  $T < S < C$  and  $T < C < S$ .

Other causal orderings, such as  $S < C < T$  and  $C < S < T$  are such that the conditional independence constraint does not lead to any simplification of the causal structure. For instance, for  $S < C < T$ , the joint distribution decomposes as  $P(S, T, C) = P(S)P(C|S)P(T|C, S)$ , and none of the terms on the right-hand side can be simplified by  $(S \perp C|T)$ . These two orderings lead to the two causal structures in Fig. 11.

Therefore, in this example, the six possible causal orderings have led to five candidates for the causal structure, depicted in Figs 7, 9, 10 and 11. However, the two



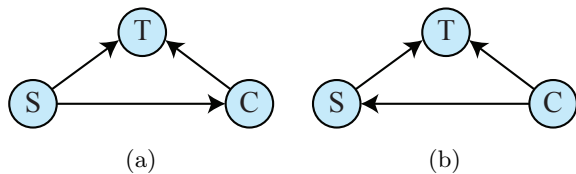


FIG. 11: DAGs that capture  $(S \perp C | T)$  for the causal orderings  $S < C < T$  (11a) and  $C < S < T$  (11b).

causal structures shown in Fig. 11 do not satisfy stability, so only the other three are viable.

Suppose finally that in addition to the information about conditional independence, one has information which rules out certain causal orderings. For instance, in the example we are considering, suppose one has the additional information that tar in the lungs always appears *after* a person has smoked, never before. It is then reasonable to rule out any causal structure that has  $T < S$ . This rules out Figs 9 and 10. At the end, the only candidate causal structure which is left is the one given in Fig. 7, which says that smoking causes tar in the lungs which causes lung cancer.

Of course, it needn't be the case that these observed variables are the only ones that are causally relevant. For instance, there might be an unobserved genetic factor which predisposes people both to smoke and to develop lung cancer. Indeed, Tobacco companies were quick to point out the possibility of explaining the observed correlation between smoking and cancer in terms of such a genetic factor. So it is useful also to have causal discovery algorithms that allow for latent variables.

Before moving on to algorithms that posit latent variables, we pause to note that the algorithm described here is proven to be correct in the sense that if there exists a set of causal structures that are minimal and faithful to the observed correlations, then the algorithm will return these structures [8].

More efficient versions of this algorithm are described elsewhere, for instance, the Inductive causation (IC) algorithm described in Pearl [1], which is equivalent to the SGS algorithm of Spirtes, Glymour and Scheines [2]. There have also been many proposals to further improve the efficiency of these algorithms (See Refs. [1] and [2] for details). These algorithms have been proven to be correct in the sense that if there exist causal models that are minimal and faithful, then the algorithms will return them.

## B. Example of causal discovery allowing for latent variables

Causal discovery in the case where one allows latent variables is more complicated. We begin by considering some of the consequences of the assumption of minimality for causal models with latent variables.

First of all, it is clear that one needn't consider any causal models wherein a latent variable *mediates* a relation between two observed variables, because the set of distributions over the observed variables that can be explained by such a model is no greater than the set that can be explained by simply postulating a direct causal influence between the observed variables. Similarly, positing a latent variable that is a common *effect* of the observed variables does not change the distributions that can be supported on the observed variables. Latent variables have nontrivial consequences for the observed distribution only when they act as *common causes* of the observed variables.

Consider the following suggestion for a causal explanation of the correlations among a set of observed variables: there are no causal influences among any of the observed variables, but there is a single latent variable that has a causal influence on each of them. By choosing the latent variable to take as many values as there are valuations of the observed variables, one can explain *any* correlation among the observed variables in this way. However, if there exists another causal model that can only reproduce a smaller set of possible correlations, while reproducing the observed correlations, then Occam's razor dictates that we should prefer the latter. Of course, one could imagine that further investigations (involving interventions, for instance) might vindicate the explanation that is less falsifiable over the one that is more falsifiable. This simply is another reminder that causal discovery algorithms are not infallible — they are heuristics for identifying the most plausible causal explanations given the evidence.

Now we come to the most subtle part of the causal discovery algorithms that posit latent variables. There is a difference between applying the criterion of minimality among a set of causal structures that are consistent with a given *distribution* over the observed variables and applying the criterion of minimality among a set of causal structures that are consistent with a given set of *conditional independence relations* over the observed variables. As we've mentioned before, the algorithms described in Refs. [1] and [2] look only at the CI relations and consequently they follow the latter course. This choice is a significant shortcoming of current causal discovery algorithms, but we will defer this criticism until the end of this section.

For the moment, we simply explain the consequences of this choice. To do so, it is useful to divide the causal structures that are consistent with a given distribution over a set of observed variables into two sorts. The first kind is such that all the latent variables it posits are common causes for at most two of the observed variables. We'll say that such a causal structure is limited to *pairwise confounding*. The other kind is unrestricted, so that more than two observed variables can be directly influenced by a single latent variable.

It is possible to show [9] that for a given set of CI relations among a set of observed variables, if a causal

model  $M$  generates those CI relations faithfully (that is, as a consequence of the causal structure, rather than the causal parameters), then there is another causal model  $M'$  that achieves the same CI relations faithfully but which is limited to pairwise confounding. The assumption of minimality makes  $M'$  preferred to  $M$ .

Therefore, if one is only applying the criterion of minimality among a set of causal structures that are consistent with the CI relations among the observed variables, then Occam's razor dictates that one need only look among causal models that are pairwise confounding. This is precisely what the standard causal discovery algorithms do. As such, one can make use of a simplified graphical language to express the set of causal structures that can be output by these algorithms. Rather than using DAGs that include both the latent and the observed variables in the causal structure, it is convenient to use a graph which only includes the observed variables as nodes but uses a larger variety of edges among these nodes to specify the causal relation that might hold among the associated variables. For instance, a double-headed arrow between variables  $X$  and  $Y$  signifies that there is a common cause of  $X$  and  $Y$  (Fig. 12). An arrow that has a circle rather than an arrowhead at one end represents either a common cause or a direct causal influence or both (Fig. 13). Finally, an undirected edge with a circle at its head and tail represents any of the five possible ways in which a pair of variables might be related (Fig. 14). In this way, a set of causal structures that include latent variables can be summarized in a single graph. Following Pearl, we call such graphs *patterns*<sup>5</sup>.

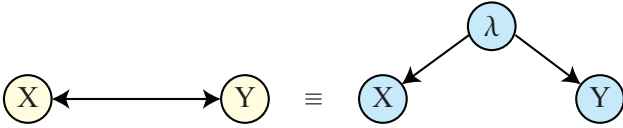


FIG. 12: The interpretation of a bidirected edge in terms of a DAG.

In order to infer which set of causal structures (including latent variables) is consistent with a given pattern, it is *not* sufficient to simply substitute for every undirected edge (or bi-directed edge or directed edge with decorated tail) all the possibilities consistent with that edge, as enumerated in Figs. 12, 13 and 14. One must eliminate some of the combinations. The definition of a *v-structure* in a DAG is a head-to-head collision of two arrows on a node such that the parents do not exert any direct causal influence on one another. The prescription

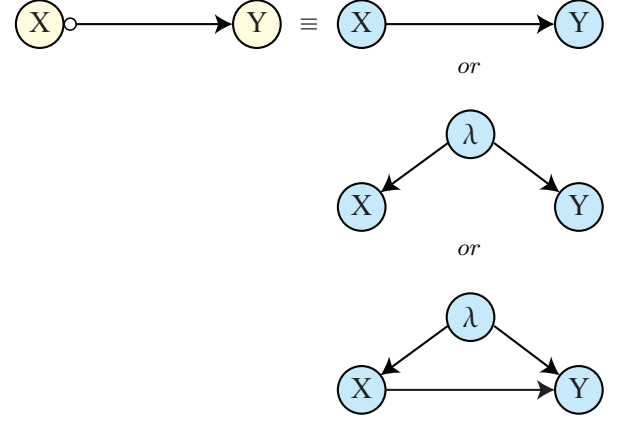


FIG. 13: The interpretation of a directed edge with a circle at its tail in terms of DAGs.

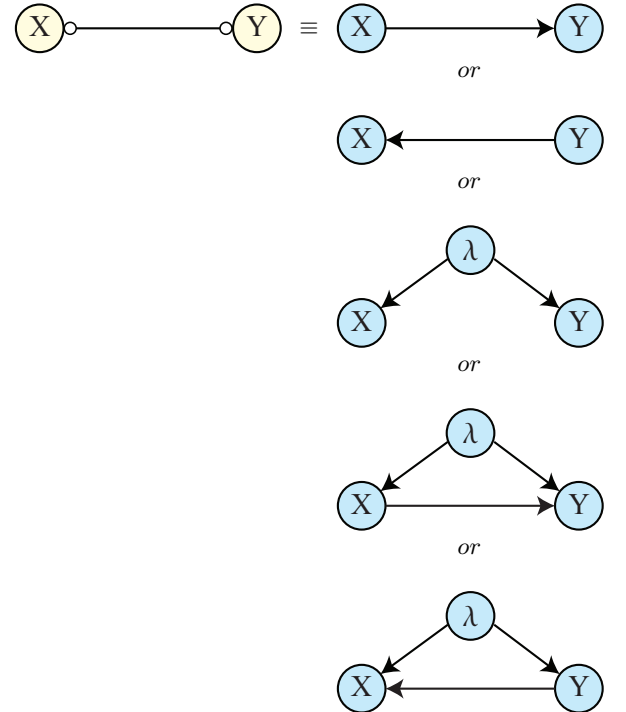


FIG. 14: The interpretation of an undirected edge with circles at head and tail in terms of DAGs.

for finding all the DAGs consistent with a pattern is to consider all the combinations of possibilities that *do not* create a new *v-structure*.

The IC\* algorithm described in Pearl [1] (which is equivalent to the Causal Inference (CI) algorithm described in SGS [2]) takes conditional independence relations as input and returns a pattern. This algorithm is correct in the sense that if there exist causal structures that are faithful to the observed CI relations, then the

<sup>5</sup>More precisely, the analogue of the particular graphs we consider here are Pearl's "marked patterns". These have also been called "partially oriented inducing path graphs" in SGS. We will follow the notational convention of SGS rather than those of Pearl when drawing such graphs.

algorithm will return the minimal structures within this set. We will not review the details of the algorithm here, but we will apply it to a simple example to get a feeling for how it works.

Consider the smoking example again, where the observed variables  $S, T$  and  $C$  are found to satisfy  $S \perp C|T$ . The pattern returned by the IC\* algorithm in this case is shown in Fig. 15.

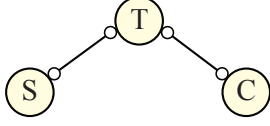


FIG. 15: Output pattern of IC\* algorithm for input  $S \perp C|T$ .

For each undirected edge in this pattern, there are five possibilities in the DAG for what connection holds between the nodes, as displayed in Fig. 14. In Fig. 16 we display all twenty-five combinations of such possibilities. We have also shaded out each of the combinations that introduces a new v-structure – these combinations are *not* candidates for the causal structure according to the IC\* algorithm. Hence, the nine causal structures that remain are the candidates returned by IC\*.

How does this answer embody the principles of causal discovery? First, the fact that the pattern admits only pairwise confounding is a consequence of the particular sort of minimality assumption going into these algorithms, as we discussed at the beginning of this section. This is the reason that we do not find in the output of the algorithm any latent variable that is a common cause of all *three* variables  $S, T$  and  $C$ .

Now consider the question of why there is neither a direct causal influence between  $S$  and  $C$  nor a latent variable that acts as a common cause for the pair. The answer is simply that if either of these sorts of influences were acting, then we would *not* find  $(S \perp C|T)$ ; learning  $S$  *would* teach us something about  $C$  even though  $T$  is known. In the context of our example, this eliminates the possibility put forward by the tobacco companies of a hypothetical genetic factor that both predisposes people to smoke and to get lung cancer.

We need not consider the cases where there is also no connection between  $S$  and  $T$  nor the cases where there is also no connection between  $T$  and  $C$  because by assumption  $(S \perp C|T)$  is the *only* CI relation and therefore  $(S \not\perp T)$  and  $(T \not\perp C)$ .

It follows that the twenty-five structures displayed in Fig. 16 are the only possibilities that remain among all possible causal structures, so to explain why the output of the algorithm is justified we need only explain why we should eliminate those that introduce a new v-structure. First note that if one conditions on a variable that is the common effect of two other variables, then we expect a dependence between those variables (for instance, in digital logic, knowing that the output of an AND gate is 0

implies that the two inputs cannot both be 1). Therefore for each causal structure that includes a v-structure on  $T$ , we would expect that conditioning on  $T$  induces a dependence between the roots of the v-structure, and because one of these roots is always correlated with  $S$  and the other with  $C$ , this would imply a dependence between  $S$  and  $C$ , contradicting the fact that  $(S \perp C|T)$ . Alternatively, we can infer that a causal structure including a v-structure on  $T$  contradicts the relation  $(S \perp C|T)$  using the d-separation criterion.

What does this imply about whether smoking causes lung cancer? Suppose that we make use of the same additional information as we considered in Sec. III A, namely, that tar in the lungs is always found to occur *after* smoking, never before. We can then eliminate all causal structures with an arrow from  $T$  to  $S$ . What remains are the three options in Fig. 17. They are: (i) smoking causes tar in the lungs which causes cancer, (ii) there is a latent variable that is a common cause of smoking and having tar in the lungs, and (iii) both mechanisms are in play. If option (ii) holds then smoking is *not* a cause of cancer and, unlike the hypothesis of a genetic factor that predisposes people both to smoke and to develop lung cancer, it *is* consistent with the observation that tar screens off smoking from cancer. Of course, this hypothesis remains implausible if one cannot identify (or imagine) any factor that screens off smoking from tar in the lungs.

We previously highlighted the fact that the causal discovery algorithms of Refs. [1] and [2] apply the principle of minimality within the set of causal structures that are consistent with the CI relations in the observed distribution, not within the set of those that are consistent with the observed distribution itself. This can be a problem because these two sets of causal structures can be different [9].

It is best to illustrate this with an example. Consider the case of a triple of observed variables,  $X, Y$  and  $Z$ . We will compare two causal models. The first posits a latent variable  $\lambda$  which has a direct causal influence on all three observed variables. The second posits three latent variables,  $\lambda, \mu$  and  $\nu$ , each of which has a direct causal influence on a distinct pair of observed variables<sup>6</sup>. The two models are illustrated in Fig. 18.

The two structures imply precisely the same set of CI relations among the observed variables, namely, the null set. However, there are distributions over the triple of observed variables that are only consistent with the first model and not the second. For instance, a joint distribution wherein the three observed variables  $X, Y$  and  $Z$  are close to perfectly correlated<sup>7</sup> cannot be generated

<sup>6</sup>This causal scenario has also been considered in the context of a discussion of nonclassical correlations in Ref. [10].

<sup>7</sup>We cannot take the case where they are perfectly correlated because we want our example to be of a distribution that is faithful to the first causal structure and perfect correlation would imply that any two variables are conditionally independent given the third.

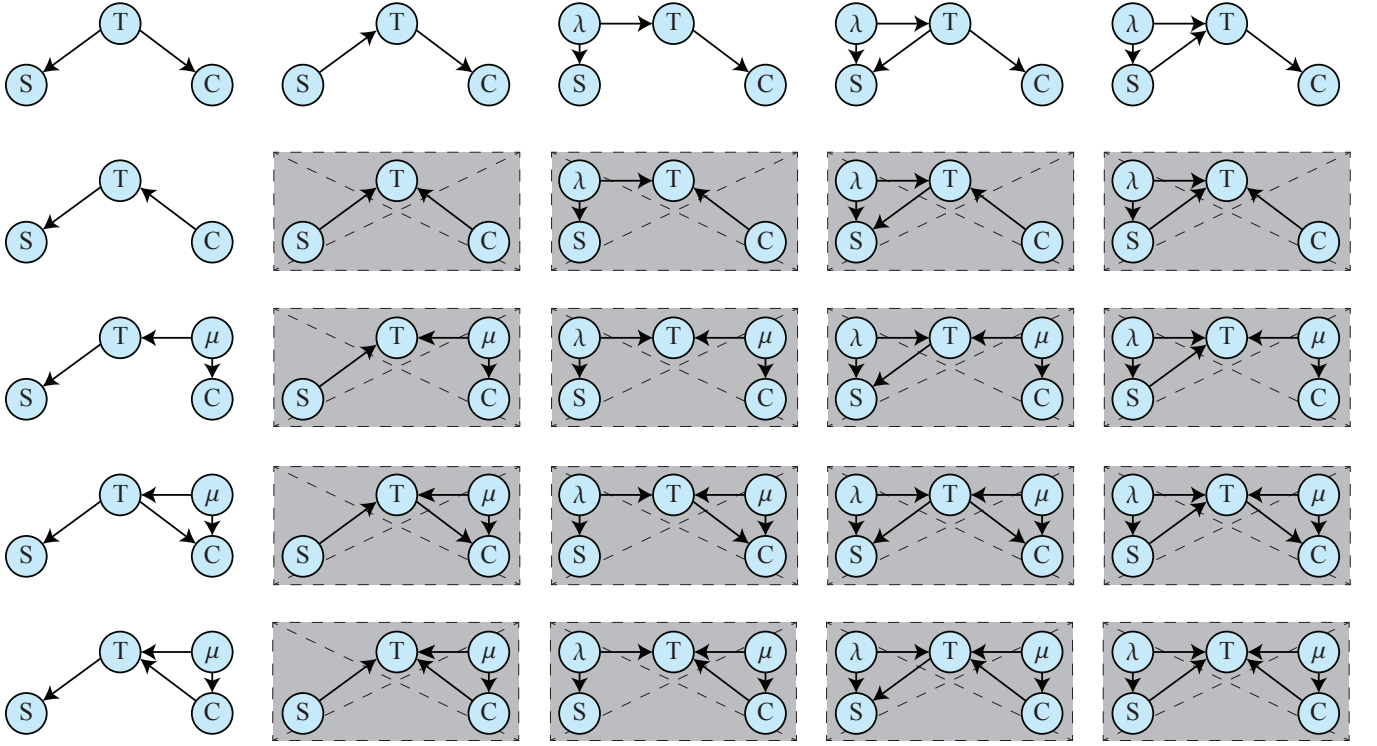


FIG. 16: The causal structures returned by the IC\* algorithm when the input is a distribution over observed variables  $S$ ,  $T$  and  $C$  with  $(S \perp C|T)$ . Those that introduce a new v-structure are shaded out.

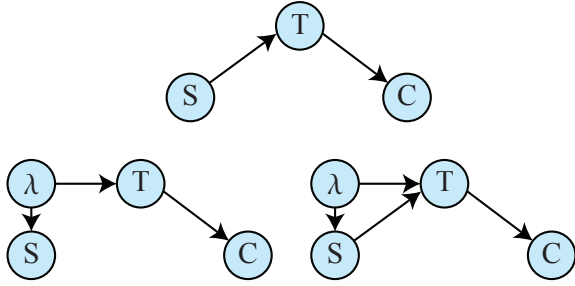


FIG. 17: The causal structures that remain if the ordering  $S < T$  is assumed.

from the second causal structure for any choice of causal parameters. Therefore, if this is the distribution one has observed, then the second causal structure is not a candidate for the underlying causal model. However, the CI relations one observes for such a distribution *are* consistent with the second causal structure. So if the input to one's causal discovery algorithm is limited to these relations, the algorithm can return a causal structure that is inconsistent with the observed distribution. Indeed, Occam's razor prefers the second structure to the first, so the causal algorithms would output a causal structure that is actually inconsistent with the observed distribution.

We will see that this sort of failure mode of the causal discovery algorithms is exactly what occurs when one

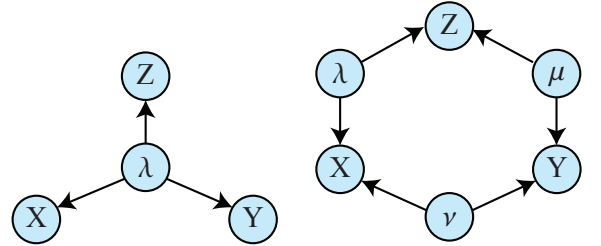


FIG. 18: Two candidate causal structures for explaining correlations between  $X$ ,  $Y$  and  $Z$  using latent variables.

applies them to correlations that violate a Bell inequality.

#### IV. APPLYING CAUSAL DISCOVERY ALGORITHMS TO QUANTUM CORRELATIONS

We now turn to the question of what these algorithms tell us about quantum correlations. We consider only Bell-type experiments involving two systems, two possible settings for each measurement and two possible outcomes for each measurement. Let  $S$  and  $T$  be the bit-valued variables that specify which measurement was performed on the left and right wings of the experiment respectively, and let  $A$  and  $B$  be the bit-valued variables that specify the outcomes of the measurements on the left and right wings respectively.

Bell's theorem derives constraints on  $P(AB|ST)$  from assumptions about the causal structure [11]. These assumptions — which Bell justified by appeal to the space-like separation of the two wings of the experiment and the impossibility of superluminal causal influences — are that  $A$  is the joint effect of the setting variable  $S$  and a common cause variable  $\lambda$ , while  $B$  is the joint effect of the setting variable  $T$  and  $\lambda$ . The causal structure corresponding to this assumption is presented in Fig. 19.

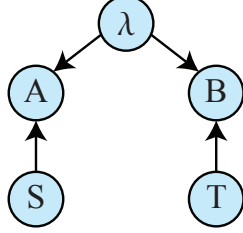


FIG. 19: The causal structure corresponding to Bell's notion of local causality.

This structure implies the following conditional independence relations,

$$(A \perp BT | S\lambda) \text{ and } (B \perp AS | T\lambda).$$

Bell called his assumption *local causality* and formalized it in terms of these conditional independences. These in turn imply that  $P(AB|ST\lambda) = P(A|S\lambda)P(B|T\lambda)$ , which is known as factorizability. From this condition, together with the assumption that there are no correlations between the settings and the hidden variables,

$$(S \perp T\lambda) \text{ and } (T \perp S\lambda),$$

one can infer that  $P(AB|ST)$  must satisfy the Bell inequalities [11, 12]. Bell's assumption about the causal structure also implies no superluminal signalling:

$$\text{No-signalling: } (A \perp T | S) \text{ and } (B \perp S | T). \quad (4.1)$$

The fact that quantum correlations can violate Bell inequalities shows that they cannot be explained using the causal structure of Fig. 19.

We will now consider the inverse problem to the one considered by Bell. Rather than attempting to infer constraints on correlations from assumptions about the causal structure, we will attempt to infer conclusions about possible causal structures from the nature of the correlations implied by quantum theory. This is the sort of problem that the causal discovery algorithms were designed to solve.

We will contrast two examples of quantum correlations: one which violates the Bell inequalities and the other which satisfies the Bell inequalities.

For the latter, we will take a version of the Einstein-Podolsky-Rosen (EPR) experiment [13] in terms of qubits

(first proposed by Bohm for spin-1/2 systems [14]). The pair are prepared in the maximally entangled state

$$|\Psi\rangle = \frac{1}{\sqrt{2}}(|+z\rangle|+z\rangle + |-z\rangle|-z\rangle) \quad (4.2)$$

where  $|\pm z\rangle$  are the eigenstates of spin along the  $\hat{z}$  axis. On each wing, the two choices of measurement are between the same pair of mutually unbiased bases, for instance, measurements of spin along the  $\hat{z}$  or  $\hat{x}$  axes, as illustrated in Fig. 20. In this case, if the same measurement is made on both wings (both  $\hat{z}$  or both  $\hat{x}$ ), one sees perfect correlation between the outcomes, while if *different* measurements are made ( $\hat{z}$  on one and  $\hat{x}$  on the other), then one sees no correlation between the outcomes. It is well known that these sorts of correlations *do not* violate any Bell inequality, which is to say that they can be explained by a locally causal model.

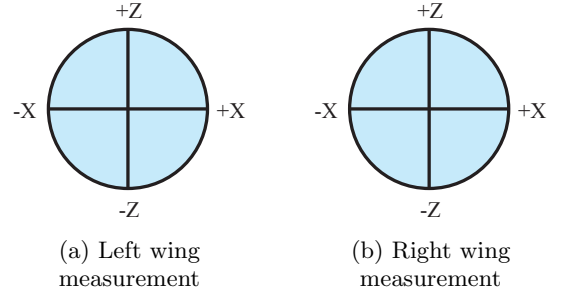


FIG. 20: Measurement axes for generating EPR correlations given the quantum state  $|\Psi\rangle$  of Eq. (4.2)

The other sort of correlation we consider will be those exhibited in the Clauser-Horne-Shimony-Holt (CHSH) experiment. We can take the pair of spins to be prepared in the same maximally entangled state  $|\Psi\rangle$  as for the EPR scenario, and the pair of measurements on the left wing to also be of spin along the  $\hat{z}$  or  $\hat{x}$  axes. However, on the right wing, the pair of possible measurements are of spin along the  $(\hat{z} + \hat{x})/\sqrt{2}$  axis or along the  $(\hat{z} - \hat{x})/\sqrt{2}$  axis, as indicated in Fig. 21. In this case, one finds that the probability of correlation for the cases  $(S, T) = (0, 0), (1, 0)$  and  $(0, 1)$  is equal to the probability of anticorrelation for the cases  $(S, T) = (1, 1)$  and has the value  $\frac{1}{2} + \frac{1}{2\sqrt{2}} \simeq 0.85$ .

The input to the standard causal discovery algorithms is limited to conditional independence relations, so we begin by computing the conditional independences that hold for the EPR and CHSH experiments. Rather than specifying an exhaustive list, we provide a generating set (the rest can be obtained by applying the semi-graphoid axioms). They are:

$$\text{EPR: } (S \perp T), (A \perp T | S), (B \perp S | T), \\ (AB \perp S), (AB \perp T).$$

$$\text{CHSH: } (S \perp T), (A \perp T | S), (B \perp S | T).$$

Consider the conditions  $(AB \perp S)$  and  $(AB \perp T)$  seen in the EPR experiment. These imply, by decomposition,



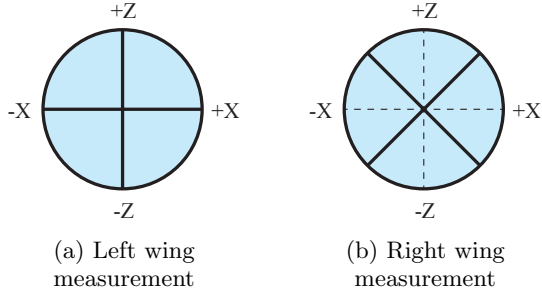


FIG. 21: Measurement axes for generating CHSH correlations given the quantum state  $|\Psi\rangle$  of Eq. (4.2)

that  $(A \perp S)$  and  $(B \perp T)$ ; the outcome on a wing is independent of the setting on that wing. While true, this independence is not representative of the causal structure. Indeed, it only holds because of the degeneracy of the Schmidt coefficients in the maximally entangled state. If we instead consider the state

$$|\Psi\rangle = \sqrt{p}|+z\rangle|+z\rangle + \sqrt{1-p}|-z\rangle|-z\rangle$$

where  $p \neq 1/2$ , then  $A \not\perp S$  and  $B \not\perp T$ . Because it is intuitively clear that the choice of measurement *does* have a causal influence on the outcome, the independences  $(A \perp S)$  and  $(B \perp T)$  are pathological in the context of the causal discovery algorithms. Given that the EPR (CHSH) experiment with a state that is close to maximally entangled still satisfies (violates) the Bell inequalities, we consider these states instead (If one likes,  $p$  may be taken to be arbitrarily close to  $1/2$ ).

We then get the following generating sets of independence relations,

$$\begin{aligned} \text{EPR:} \quad & (S \perp T), (A \perp T|S), (B \perp S|T), \\ \text{CHSH:} \quad & (S \perp T), (A \perp T|S), (B \perp S|T), \end{aligned}$$

where  $(S \perp T)$  asserts the independence of the settings, and  $(A \perp T|S)$  and  $(B \perp S|T)$  are the no-signalling conditions. The critical point is that the set of independences are the *same* for the EPR and the CHSH experiments. Because the input to the causal discovery algorithms that we consider is limited to conditional independence relations, it follows that whatever causal conclusions these algorithms draw, they will draw the *same* causal conclusions about the EPR experiment as they do about the CHSH experiment. And yet, from the fact that the EPR correlations satisfy the Bell inequalities, we know that they *can* be explained by local causes while from the fact that the CHSH correlations violate a Bell inequality, we know that they cannot be so explained.

So the conclusion is that *standard causal discovery algorithms (based on conditional independences) cannot possibly do justice to Bell's theorem*. Independences simply do not provide enough information. One needs a causal discovery algorithm that looks at the strength of correlations to reproduce the conclusions of Bell's theorem.

Despite the inability of the standard causal discovery algorithms to distinguish correlations that violate the Bell inequalities from those that satisfy them, it is nonetheless interesting to see what happens when one applies the algorithms to the set of independences we found for the EPR and CHSH experiments. We will refer to these as *nontrivial no-signalling correlations* (they are deemed nontrivial because they predict correlation between the outcomes for some choices of the settings).

In applying the causal discovery algorithms, we will assume that the setting variable on one wing is a cause of the outcome variable on that wing, that is, we assume that  $S$  is a cause of  $A$  and that  $T$  is a cause of  $B$ . This is presumably uncontroversial. The assumption that there are no causal cycles then implies that there can be no causal influence from  $A$  to  $S$ , nor from  $B$  to  $T$ . Nonetheless, we are still permitting influences from the outcome on one wing to the setting on the other, although, as we will see, the causal discovery algorithms will rule against such influences.

### A. No latent variables

It is instructive to consider the causal structure that arises for a single representative causal ordering of the variables. We take  $S < T < A < B$ . Then, the most general causal structure is illustrated in Fig. 22. Hence the most general joint distribution for this ordering is of the form

$$P(S, T, A, B) = P(S)P(T|S)P(A|S, T)P(B|S, T, A).$$

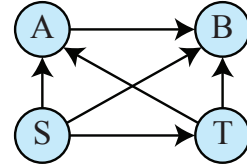


FIG. 22: The most general causal structure for the causal ordering  $S < T < A < B$ , assuming no hidden variables.

The independence  $(S \perp T)$  implies that  $P(T|S) = P(T)$ , and the independence  $(A \perp T|S)$  implies that  $P(A|S, T) = P(A|S)$ . The independence  $(B \perp S|T)$  has no nontrivial implications for this causal ordering, hence the term  $P(B|S, T, A)$  cannot be simplified. From these CI relation it follows that the joint distribution can be written as

$$P(S, T, A, B) = P(S)P(T)P(A|S)P(B|S, T, A),$$

which corresponds to the causal structure in Fig. 23a. If we change the ordering of variables so that  $B$  precedes  $A$ , then by a similar argument, we obtain the causal structure in Fig. 23b. For every other possible causal ordering



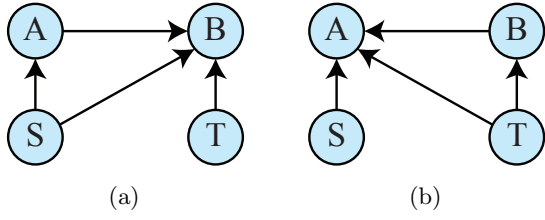


FIG. 23: Possible causal structures for no-signalling correlations, assuming no hidden variables, for causal orderings  $S < T < A < B$ ,  $T < S < A < B$ ,  $S < A < T < B$  (23a) and  $S < T < B < A$ ,  $T < S < B < A$ ,  $T < B < S < A$  (23b).

consistent with our assumption that  $S < A$  and  $T < B$ , we also obtain one of the causal structures of Fig. 23.

Consider the causal structure in Fig. 23a. Although it faithfully captures  $(S \perp T)$  and  $(A \perp T|S)$ , it does not faithfully capture  $(B \perp S|T)$ . The only way to explain the independence  $(B \perp S|T)$  within this causal model is by fine-tuning of the causal parameters in the model, for instance, if the parameters defining  $P(B|S, T, A)$  are not independent of those defining  $P(A|S)$ . A similar problem arises for the causal structure in Fig. 23b. It follows that in the case of no latent variables no causal structure can satisfy minimality and faithfulness for the conditional independences of nontrivial no-signalling correlations.

Note that if, instead of applying the Wermuth-Lauritzen algorithm to the nontrivial no-signalling correlations, one applies the IC algorithm [1], equivalently the SGS algorithm [2], one finds that it returns a graph that is not a valid causal structure, signalling a failure of the algorithm. This is what one would expect given that the algorithm only promises to return a valid causal structure if there exists one that is minimal and faithful to the correlations, and in this case, there is not.

There is an interesting lesson here for the foundations of quantum theory. Long before Bell's work, Einstein had pointed out that if one did *not* assume hidden variables, then one could only explain the EPR correlations by positing superluminal causes. This argument was made in his comments at the 1927 Solvay conference [15] (See Refs. [16] and [3] for more concerning Einstein's arguments on completeness and locality.) One can easily cast Einstein's argument into the mold of causal structures as follows. If we allow the quantum state  $\psi$ , considered as a classical variable, as the only common cause, then the assumption of no superluminal causes implies that  $P(A, B|S, T, \psi) = P(A|S, \psi)P(B|T, \psi)$ , and given that  $\psi$  is fixed in the experiment (it is a variable which only takes one possible value), this implies that  $A$  and  $B$  should be uncorrelated, in contradiction with the EPR correlations.

But what the result of our analysis shows is that Einstein failed to explicitly note another mysterious feature of the EPR correlations, namely, that *even if* one was willing to countenance superluminal causes in an attempt to explain the EPR correlations without recourse to hid-

den variables, ensuring that these superluminal causes cannot be used to send superluminal signals implies that there must be fine-tuning in the underlying causal model.

## B. Latent variables allowed

If one simply inputs the independences of nontrivial no-signalling correlations into the IC\* algorithm of Ref. [1], one obtains the causal diagram illustrated in Fig. 24 as output.

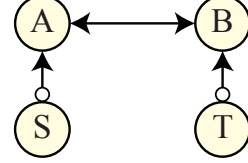


FIG. 24: The output pattern of the IC\* algorithm when applied to nontrivial no-signalling correlations.

Recall that the arrows with an empty circle at their tail imply that one can have either a direct causal link or a common cause. If one believes that the settings at each wing are freely chosen, then one is inclined to think that either the setting variables  $S$  and  $T$  should be direct causes of  $A$  and  $B$  respectively, or that if they are not, then it is the common cause for  $A$  and  $S$  and the common cause for  $B$  and  $T$  that is freely chosen. In this case, we could lump the common causes into the definition of the setting variables without loss of generality.

Besides this caveat about the causal relation between  $S$  and  $A$  and between  $T$  and  $B$ , the causal structures consistent with the pattern that the IC\* algorithm has returned are precisely those that capture Bell's notion of local causality, illustrated in Fig. 19.

Recall that the CHSH correlations are an instance of nontrivial no-signalling correlations that *violate* the Bell inequalities. Therefore, the IC\* algorithm is claiming that Bell-inequality-violating correlations can be explained by a locally causal model. However from Bell's theorem we know that this claim is mistaken. Therefore, the IC\* algorithm comes to a causal conclusion that is incorrect.

Of course, we already pointed out in Sec. IV, that the input of the IC\* algorithm cannot distinguish Bell-inequality-violating from Bell-inequality-satisfying correlations. The independences we have fed into the algorithm also hold for EPR correlations. Consequently, had it returned the conclusion that the nontrivial no-signalling correlations *cannot* be explained by the causal structure of Fig. 19, it would have also reached an incorrect causal conclusion because the EPR correlations *can* be explained by such a model.

So we reiterate our conclusion from Sec. IV, that causal discovery algorithms which look only at independences are inadequate to the task of establishing whether or not

correlations can be explained by a locally causal model. We require better algorithms that also take into account the strengths of the correlations.

From our brief discussion in Sec. III B of the shortcomings of causal discovery algorithms with latent variables we can also see *why* the algorithms have reached an incorrect conclusion. The problem is that a causal structure with latent variables that reproduces the CI relations of a given distribution might not be capable of reproducing the distribution itself. In particular, the causal structure of Fig. 24 reproduces the CI relations of the distribution  $P(A, B, S, T)$  defined by the CHSH experiment, but it cannot reproduce the distribution itself.

### C. Some proposed causal explanations of quantum correlations

We now apply the ideas behind causal discovery algorithms to a few of the existing proposals for providing a causal explanation of Bell-inequality-violating correlations. We consider three: superluminal causation, superdeterminism, and retrocausation.

#### 1. Superluminal causation

One option for explaining Bell correlations causally is to assume that there are some superluminal causes, for instance, a causal influence from the outcome on one wing to the outcome on the other, or from the setting on one wing to the outcome on the other, or both. In the most general case one allows hidden variables that can causally influence the measurement outcomes. The possibilities are illustrated in Fig. 25.

But the same problem arises for these sorts of causal explanations of Bell-inequality violations as arise for the causal explanations without hidden variables that were discussed in Sec. IV A. Given the superluminal causes from one wing to the other, the only way to explain the lack of superluminal signals is through a fine-tuning of the causal parameters.

For instance, in Fig. 25c, the correlations set up between  $S$  and  $B$  along the direct causal path could cancel with those set up by the causal path through  $A$ . (The path through  $\lambda$  cannot set up correlations between  $S$  and  $B$  because there is a collider on  $A$  in this path and we are not conditioning on  $A$ .) Such a cancellation requires fine-tuning of the parameters of the model.

To salvage no-signalling for the causal structure of Fig. 25a, we need a different sort of fine-tuning (a similar sort of fine-tuning mechanism can also be used for the causal structure of Fig. 25b). For instance, it could be that  $\lambda = (\lambda_1, \lambda_2)$  where  $\lambda_1$  is a binary variable that is uniformly distributed and that  $Y$  is a function of  $S \oplus \lambda_1$ ,  $T$  and  $\lambda_2$ . In this case, we can ensure that  $(Y \perp S|T)$  by virtue of the special distribution on  $\lambda_1$ , which is a kind of fine-tuning.

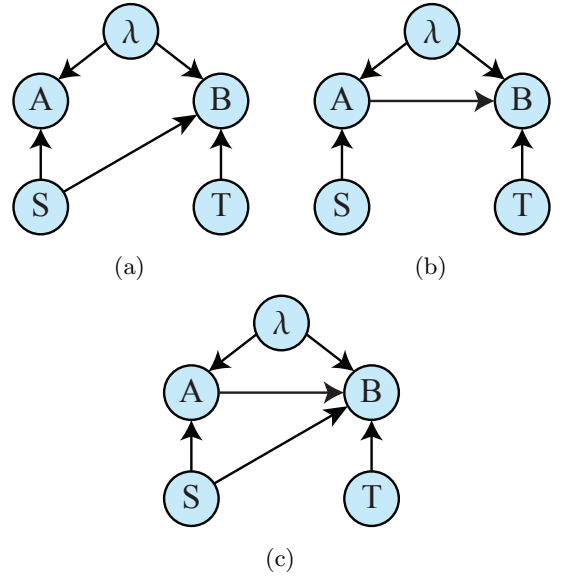


FIG. 25: Examples of causal structures that posit superluminal causal influences to explain Bell correlations.

Note that this is precisely the sort of causal structure that is assumed in the Toner and Bacon model [17], where Bell-inequality violations are simulated by classical communication<sup>8</sup>. This model also involves fine-tuning insofar as signalling is prohibited only for a special distribution over the shared random variables posited by the model.

The deBroglie-Bohm interpretation is a prominent example of a model that seeks to provide a causal explanation of Bell correlations using superluminal causal influences. Consider the deBroglie-Bohm interpretation of a relativistic theory such as the model of QED provided by Struyve and Westman [18], or else of a nonrelativistic theory wherein the interaction Hamiltonians are such that there is a maximum speed at which signals can propagate. In both cases, it is presumed that there is a preferred rest frame that is hidden at the operational level. In a Bell experiment, if the measurement on the left wing occurs prior to the measurement on the right wing relative to the preferred rest frame, then there is a superluminal causal influence from the setting on the left wing to the outcome on the right wing, mediated by the quantum state, which is considered to be a part of the ontology of the theory [19]. (Note that no causal influence from the outcome of the first experiment to the outcome of the second is required because the outcomes are deterministic functions of the Bohmian configuration and the wavefunction.) It follows from our analysis that the parameters in the causal model posited by the deBroglie-Bohm inter-

<sup>8</sup>This model works even when the measurement setting for each qubit is chosen arbitrarily, rather than being limited to the two settings of the CHSH experiment.

pretation must be fine-tuned in order to explain the lack of superluminal signalling.

Valentini's version of the deBroglie-Bohm interpretation makes this fact particularly clear. In Refs. [20, 21] he has noted that the wavefunction plays a dual role in the deBroglie-Bohm interpretation. On the one hand, it is part of the ontology, a pilot wave that dictates the dynamics of the system's configuration (the positions of the particles in the nonrelativistic theory). On the other hand, the wavefunction has a statistical character, specifying the distribution over the system's configurations. In order to eliminate this dual role, Valentini suggests that the wavefunction is only a pilot wave and that *any* distribution over the configurations should be allowed as the initial condition. It is argued that one can still recover the standard distribution of configurations on a coarse-grained scale as a result of dynamical evolution [22]. Within this approach, the no-signalling constraint is a feature of a special equilibrium distribution. The tension between Bell inequality violations and no-signalling is resolved by abandoning the latter as a fundamental feature of the world and asserting that it only holds as a *contingent* feature. The fine-tuning is explained as the consequence of equilibration. (It has also been noted in the causal model literature that equilibration phenomena might account for fine-tuning of causal parameters [23].) Conversely, the version of the deBroglie-Bohm interpretation espoused by Dürr, Goldstein and Zhang [24] – which takes no-signalling to be a non-contingent feature of the theory – does not seek to provide a dynamical explanation of the fine-tuning. Consequently, it seems fair to say that the fine-tuning required by the deBroglie-Bohm interpretation is less objectionable in Valentini's version of the theory.

## 2. Superdeterminism

Another option for a causal explanation of quantum correlations is to posit that the settings are not free but are causally influenced by other variables.

For instance, the hidden variable  $\lambda$  (which correlates the outcomes) might causally influence one or both of the setting variables, as illustrated in Figs. 26a and 26b. Alternatively, one can posit the existence of a second hidden variable  $\mu$  that is a common cause for the setting on one wing and the outcome on the other wing, as illustrated in Fig. 26c. More complicated possibilities would have  $\mu$  as a common cause of a subset of three of the settings and outcomes. Note that the possibility of a latent variable that is a common cause of  $\lambda$  and one or both settings has not been excluded; it is incorporated into the first case. This is because any such variable could just be absorbed into the definition of  $\lambda$  without loss of generality. The scenario in Fig. 26c could also be considered a special case of the one in Fig. 26a, if we include  $\mu$  into the definition of  $\lambda$ . Nonetheless, it is useful to separate out this second case because it posits that the common cause of

$A$  and  $B$  is not correlated with the common cause of  $S$  and  $T$ .

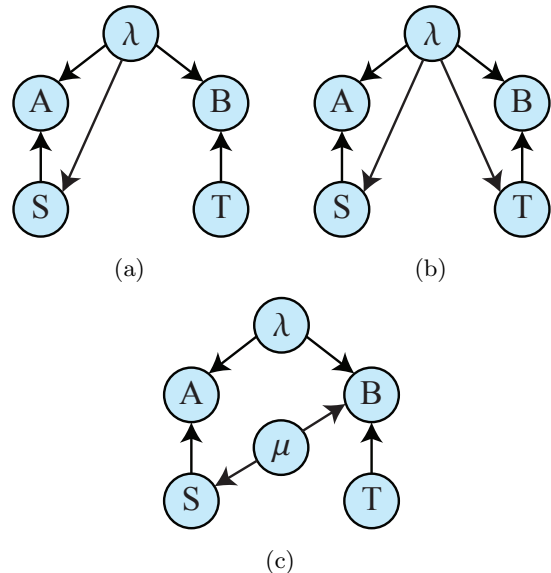


FIG. 26: Some causal structures that exploit the superdeterminism loophole to explain Bell correlations.

All of the causal influences posited in such models can be taken to be subluminal. However, such explanations of the Bell correlations are clearly in conflict with the notion that the settings can be freely chosen by the experimenter. To assert one of these causal structures as a way to resolve the mystery of Bell's theorem is an instance of what is commonly known as the “superdeterminism” loophole. But, just as with positing superluminal causal influences, these causal structures are not faithful to the observed correlations because one or more of the observed CI relations –  $S \perp T$  (independence of settings),  $(A \perp T|S)$  (no-signalling from left to right) and  $(B \perp S|T)$  (no signalling from right to left) – can only be satisfied by fine-tuning of the parameters in the causal model. This is a novel sort of objection against the notion of a superdeterministic explanation of Bell-inequality-violations, independent of an appeal to free will.

It is worth devoting a few words to the sort of fine-tuning that is required. First note that in the context of abandoning the assumption of free will, the no-signalling constraint must be reinterpreted as an observed statistical independence, rather than a statement about the consequences of an intervention on a setting variable. Of course, this statistical independence is still observed and therefore must still be reproduced by the causal model. In the causal structure of Fig. 26a, if we define  $\lambda^*$  to be that part of  $\lambda$  that is correlated nontrivially with  $S$ , then we require that  $\lambda^* \perp B$  despite the arrow from  $\lambda$  to  $B$ . We can still do justice to the Bell correlations by having  $\lambda^*$  correlated with only the *parity* of  $A$  and  $B$ , while remaining uncorrelated with  $B$ . This is an instance of

fine-tuning.

Similar fine-tuning tricks can be used to ensure that  $(B \perp S|T)$  in the causal structures of Figs. 26b and 26c.

### 3. Retrocausation

“Retrocausation” refers to the possibility of causal influences that act in a direction contrary to the standard arrow of time. It has been proposed as a means of resolving the mystery of Bell-inequality violations [25–29] by purportedly saving the relativistic structure of the theory: rather than having causal influences propagating outside the light cone, they propagate *within* the light cone although possibly within the *backward* light cone.

It is useful to distinguish two approaches to retrocausal explanations of Bell correlations: those that add cycles to the causal structure and those that do not. Given that the former take us outside the framework of directed *acyclic* graphs, we will confine our attention to acyclic retrocausal models.

Price has described the idea of a retrocausal model of Bell inequality violations in Ref. [30]. It is not completely clear whether he has in mind a model that posits cycles or not. However, he does argue that one way to generate a retrocausal model is to start with a superdeterministic model and to simply reverse the causal arrows that lead into the settings. For the examples of superdeterminism we have considered, such reversals lead to acyclic retrocausal models. For instance, if one starts with the superdeterministic causal structure of Fig. 26a and reverses the  $\lambda \rightarrow S$  arrow, one obtains the causal structure of Fig. 27a, where setting  $S$  is a cause of the hidden variable  $\lambda$ . If one assumes that  $S$  is chosen freely at a time to the future of when  $\lambda$  is set, then this model is clearly retrocausal.

Alternatively, consider taking the superdeterministic model of Fig. 26c and reversing the  $\mu \rightarrow S$  arrow, to obtain the causal structure of Fig. 27c. If  $\mu$  were presumed to be space-like separated from both  $S$  and  $B$ , it would simply mediate a superluminal causal influence from  $S$  to  $B$ . However, if one posits that  $\mu$  is in the common future of  $S$  and  $B$ , then we can imagine that there is a causal influence from  $S$  to  $\mu$  that is subluminal, and one from  $\mu$  to  $B$  that is retrocausal. Alternatively, if one posits that  $\mu$  is in the common past of  $S$  and  $B$ , then the causal influence from  $S$  to  $\mu$  must be assumed to be retrocausal.

Note that if one views spatio-temporal relations as supervening upon causal relations, rather than vice-versa, then there is no freedom to specify the spatio-temporal location of  $\mu$  and the distinction drawn above is not meaningful. Even if one takes spatio-temporal notions to be primary, the fact that the location of  $\mu$  seems to be mere window-dressing in the context of a causal explanation of Bell-inequality violations undermines the distinction between retrocausation and superluminal causation.

Fine-tuning is just as necessary within the retrocausal

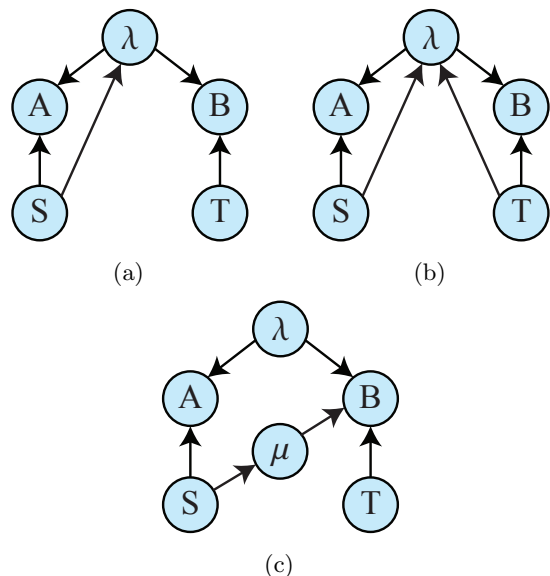


FIG. 27: Causal structures that exploit the retrocausation loophole to explain Bell correlations.

explanations as it was in the ones that posited superluminal influences or superdeterminism. Without it, one would obtain a correlation between  $S$  and  $B$ , in contradiction with their observed statistical independence. Indeed, if these causal structures could be supplemented with arbitrary causal parameters, then one could use the causal chain of influence that extends from  $S$  to  $B$  to send a signal.

## V. CONCLUSIONS

Our two main conclusions are as follows. First, causal discovery algorithms that appeal only to conditional independences among observed variables cannot distinguish between Bell-inequality-violating and Bell-inequality-satisfying correlations. Better algorithms which look to the strength of correlations are needed to do justice to Bell’s theorem.

Second, and more importantly, we have shown that any causal model which can reproduce Bell-inequality violations while respecting the observed independences—the marginal independence of the measurement settings and the no-signalling condition—will necessarily violate a principle that is at the core of all the best causal discovery algorithms, namely, that observed independences should not be explained by fine-tuning of the causal parameters in the model. This is true for all explanatory strategies that fit within the framework of directed acyclic graphs supplemented with conditional probabilities, including models that posit superluminal causes, models that exploit the superdeterminism loophole, and models that posit retrocausation while avoiding causal cycles.

The topic of causal discovery is still relatively young.

The best algorithms available today are not likely to be the final story. Indeed, our analysis suggests that the tools that have been developed in the literature on the foundations of quantum theory for assessing the possibility of *local* explanations of correlations may well be important for developing causal discovery algorithms. If one could deliver on this promise, then it would be an interesting example of the field of quantum foundations having applications in other fields, such as statistics and machine learning, and via these, in medicine, genetics, economics and other disciplines wherein causal discovery plays a prominent role.

Conversely, it is our view that there is a great deal more insight to be gained about the foundations of quantum theory from the literature on causal models and causal discovery algorithms. We consider a few possible directions of research along these lines.

As mentioned previously, defining causality in a manner that does not make reference to temporal ordering provides a language by which one could hope to describe a fundamental theory wherein spatio-temporal notions are emergent and notions of causal structure are primitive. In such a theory, it would not be the case that a cause was defined to be prior in time to its effects, but rather the notion of the temporal order of two events would be defined in terms of whether one event was a potential cause of the other. Consequently, the framework for causal inference provides a natural arena in which to pursue the idea that space-time is emergent, a notion that is popular in attempts to unify general relativity with quantum theory [31, 32].

There are a number of results in the quantum foundations literature that have the following form: make some assumptions about the causal structure and derive inequalities on the correlations that can be obtained from these classically. Svetlichny’s inequalities are an example of this [33], wherein one considers a triple of measurements at space-like separation and one allows a mixture of causal structures wherein superluminal influences can propagate between any two of the wings of the experiment. The topic has been studied in Refs. [34–36]. Fritz has also recently derived inequalities on classical correlations for some causal structures that do not correspond to the standard Bell scenario [10]. Such results are examples of a general approach to correlations that has been developed in the causal model literature. For instance, In Pearl’s book (Sec. 8.4), inequalities on correlations are derived from assumptions about the causal structure in a section considering noncompliance in drug trials. Pearl points out the similarity between these “instrumental” inequalities and the Bell inequalities, and adds: “The instrumental inequality can, in a sense, be viewed as a generalization of Bell’s inequality for cases where direct causal connection is permitted to operate between the correlated observables,  $X$  and  $Y$ .” It will be interesting to see how many results in the quantum foundations literature can be considered to be instances of such generalized inequalities.

Finally, by exploiting a quantum analogue of conditional probability proposed by Leifer [37] and developed by Leifer and Spekkens [38, 39] and an associated quantum analogue of conditional independence (see Leifer and Poulin [40], for instance), one can hope to explore a generalization of the notion of causal model to a *quantum causal model*. A quantum causal model is naturally defined as a quantum causal structure, which is a directed acyclic graph wherein each node is a quantum system, and a set of quantum causal parameters, which constitute a set of conditional quantum states (the quantum analogue of conditional probability) for every node given its causal parents. Insofar as one can accommodate classical variables as special cases of quantum systems (corresponding to commuting algebras), one can describe correlations among settings and outcomes within quantum causal models.

Quantum causal models make similar assumptions about the possibilities for causal structure as do classical causal models (no cycles for instance), and they make similar assumptions about the consequences of causal structure for statistical independences, but they replace the formalism of classical probability theory with a non-commutative generalization thereof. If one can make the case that the formalism of quantum causal models is not just a mathematical artifice but can be given a sensible interpretation as a form of causal explanation, then such models can provide a causal explanation of Bell-inequality violations without requiring fine-tuning.

Note, however, that if the conditional probabilities that appear in classical causal models are interpreted as degrees of belief – and we take this to be the most sensible interpretation – then the transition from classical causal models to quantum causal models involves not only a modification to physics, but a modification to the rules of inference. In this view, the correct theory of inference is not *a priori* but empirical. Nonetheless, one cannot simply declare *by fiat* that some formulation of quantum theory is a theory of inference. One must justify this claim. At a minimum, one must determine how standard concepts in a theory of inference generalize to the quantum domain. One could also reconsider the various proposals for axiomatic derivations of classical probability theory, for instance, that of Cox [41] or that of de Finetti [42], to see whether a reasonable modification of the axioms yields a quantum theory of inference<sup>9</sup>. Ideally, one would show that if quantum causal models imply a modification to both our physics and to our theory of inference, then these modifications are not independent. After all, the physics determines the precise manner in which an agent can gather information about the world and in turn act upon it and so the physics should deter-

<sup>9</sup>Fuchs and Schack have also suggested that parts of quantum theory can be derived by an appeal to dutch-book coherence following de Finetti [43].



mine what is the most adaptive theory of inference for an agent. It is in this sense that the project of defining quantum causal models is not yet complete and only with such a completion in hand can one really say that a causal explanation of the Bell correlations without recourse to fine-tuning has been achieved.

## VI. ACKNOWLEDGEMENTS

RWS thanks Matthew Leifer for having introduced him to causal networks and for discussions on the ideas in this

article. We would also like to thank Jonathan Barrett, Cozmin Ududec and especially Dominik Janzing for helpful discussions, and Howard Wiseman for comments on a draft of the article. Part of this work was completed while CJW was a student in the Perimeter Scholars International program. Research at Perimeter Institute is supported in part by the Government of Canada through NSERC and by the Province of Ontario through MRI. CJW is supported by the Canadian Excellence Research Chairs (CERC) Program and the Canadian Institute for Advanced Research (CIFAR).

- 
- [1] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2nd edition, 2009.
  - [2] P. Spirtes, C. N. Glymour, and R. Scheines. *Causation, Prediction, and Search*. The MIT Press, 2nd edition, 2001.
  - [3] N. Harrigan and R. W. Spekkens. Einstein, Incompleteness, and the Epistemic View of Quantum States. *Foundations of Physics*, 40(2), 2010.
  - [4] D. Janzing and J. Lemeire. Replacing causal faithfulness with algorithmic independence conditions. *To appear in Minds and Machines*, 2010.
  - [5] T. Norsen. Against realism. *Foundations of Physics*, 37(3):311–340, 2007.
  - [6] C. N. Glymour. Markov properties and quantum experiments. In *Physical Theory and its Interpretation*, pages 117–126. Springer Netherlands, 2006.
  - [7] H. Beebe. *The Oxford handbook of causation*. Oxford University Press, 2009.
  - [8] N. Wermuth and S. L. Lauritzen. Graphical and recursive models for contingency tables. *Biometrika*, 70(3), 1983.
  - [9] T. S. Verma. Graphical aspects of causal models. Technical Report R-191, Computer Science Department, University of California, Los Angeles, 1993.
  - [10] T. Fritz. Hidden bayesian networks: Inferring physical theories from correlations and causal structure. *arXiv:1206.5115v1 [quant-ph]*, 2012.
  - [11] J. S. Bell. On the einstein-podolsky-rosen paradox. *Physics*, 1(3), 1964.
  - [12] J. F. Clauser, M. A. Horne, A. Shimony, and R. A. Holt. Proposed experiment to test local hidden-variable theories. *Physical Review Letters*, 23:880–884, 1969.
  - [13] A. Einstein, B. Podolsky, and N. Rosen. Can Quantum Mechanical Description of Physical Reality be Considered Complete? *Physical Review*, 47, 1935.
  - [14] D. Bohm. *Quantum Theory*. Dover Publications, 1989.
  - [15] G. Bacciagaluppi and A. Valentini. *Quantum Theory at the Crossroads: Reconsidering the 1927 Solvay Conference*. Cambridge University Press, 2009.
  - [16] T. Norsen. J.s. bell’s concept of local causality. *arXiv:0707.0401v2 [quant-ph]*, 2010.
  - [17] B. F. Toner and D. Bacon. Communication cost of simulating bell correlations. *Phys. Rev. Lett.*, 91:187904, Oct 2003.
  - [18] W. Struyve and H. Westman. A minimalist pilot-wave model for quantum electrodynamics. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Science*, 463(2088), 2007.
  - [19] D. Bohm and B. J. Hiley. *The Undivided Universe*. Routledge, 1993.
  - [20] A. Valentini. Signal-locality, uncertainty, and the sub-quantum h-theorem. ii. *Physics Letters A*, 158(1-2), 1991.
  - [21] A. Valentini. Signal-locality, uncertainty, and the sub-quantum h-theorem. i. *Physics Letters A*, 156(1-2), 1991.
  - [22] A. Valentini and H. Westman. Dynamical origin of quantum probabilities. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Science*, 461(2053), 2005.
  - [23] D. Dash. Restructuring dynamic causal systems in equilibrium. In *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics (AISTats 2005)*, pages 81–88, 2005.
  - [24] D. Dürr, S. Goldstein, and N. Zanghi. Quantum equilibrium and the origin of absolute uncertainty. *Journal of Statistical Physics*, 67, 1992.
  - [25] O. Costa de Beauregard. Mécanique quantique. *Comptes Rendus Académie des Sciences*, 236(1632), 1953.
  - [26] J. G. Cramer. The transactional interpretation of quantum mechanics. *Reviews of Modern Physics*, 58, 1986.
  - [27] R. Sutherland. Causally symmetric bohm model. *Studies In History and Philosophy of Modern Physics*, 39, 2008.
  - [28] H. Price. Toy models for retrocausality. *Studies In History and Philosophy of Science Part B: Studies In History and Philosophy of Modern Physics*, 39(4), 2008.
  - [29] K. B. Wharton. A novel interpretation of the klein-gordon equation. *Foundations of Physics*, 40, 2010.
  - [30] H. Price. *Time’s arrow & Archimedes’ point: new directions for the physics of time*. Oxford University Press, USA, 1997.
  - [31] F. Markopoulou. New directions in background independent quantum gravity. *Arxiv preprint gr-qc/0703097*, 2007.
  - [32] F. Dowker. Causal sets and the deep structure of space-time. *100 Years of Relativity, Space-Time Structure: Einstein and Beyond*, pages 445–464, 2005.
  - [33] G. Svetlichny. Distinguishing three-body from two-body nonseparability by a bell-type inequality. *Phys. Rev. D*, 35(10), 1987.
  - [34] P. Mitchell, S. Popescu, and D. Roberts. Conditions for the confirmation of three-particle nonlocality. *Phys. Rev. A*, 70(6), 2004.
  - [35] R. Gallego, L. E. Würflinger, A. Acín, and



- M. Navascués. An operational framework for non-locality. *arXiv:1112.2647v1 [quant-ph]*, 2011.
- [36] J. Barrett, S. Pironio, J. Bancal, and N. Gisin. The definition of multipartite nonlocality. *arXiv:1112.2626v1 [quant-ph]*, 2011.
- [37] M. S. Leifer. Conditional Density Operators and the Subjectivity of Quantum Operations. *arXiv:quant-ph/0611233v1*, 2006.
- [38] M. S. Leifer and R. W. Spekkens. Formulating quantum theory as a causally neutral theory of bayesian inference. *arXiv:1107.5849 [quant-ph]*, 2011.
- [39] M. S. Leifer and R. W. Spekkens. A bayesian approach to compatibility, improvement, and pooling of quantum states. *arXiv:1110.1085v1 [quant-ph]*, 2011.
- [40] M. S. Leifer and D. Poulin. Quantum graphical models and belief propagation. *Annals of Physics*, 323(8), 2008.
- [41] R.T. Cox. Probability, frequency and reasonable expectation. *American journal of physics*, 14(1):1–13, 1946.
- [42] B. De Finetti. Foresight: Its logical laws, its subjective sources. *Breakthroughs in Statistics*, 1:134x174, 1937.
- [43] C.A. Fuchs and R. Schack. Quantum-bayesian coherence. *Arxiv preprint arXiv:0906.2187*, 2009.

## Appendix A: $d$ -separation

Conditional independence relations are captured in directed acyclic graphs by the notion of distance-separation or  **$d$ -separation**. First let us introduce the basic elements of which a DAG may be composed of; these are *colliders*, *forks*, and *chains*; which for three variables  $A, B, C$  are illustrated in Fig. 28.

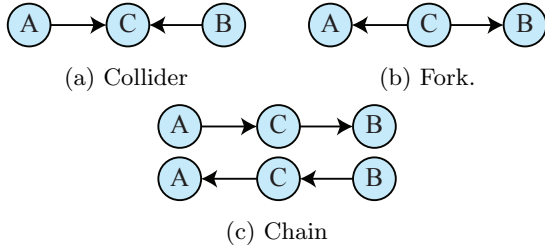


FIG. 28: Basic structures found in DAGs.

Given a DAG  $G$ , a path between two vertices  $X$  and  $Y$  is any set of edges and vertices which connects  $X$  and  $Y$ , regardless of the direction of the edges. We say that a path between  $X$  and  $Y$  is *blocked* by a set of vertices  $\mathbf{Z}$  if at least one of the following conditions holds

1. The path contains a chain (Fig. 28c), or a fork (Fig. 28b) such that  $C$  is in  $\mathbf{Z}$ .
2. The path contains a collider (Fig. 28a) such that  $C$  is not in  $\mathbf{Z}$  and no descendant of  $C$  is in  $\mathbf{Z}$ .

We then have the following definition of  $d$ -separation:

**Definition 2 ( $d$ -separation)** *Given a DAG  $G$  with vertices  $\mathbf{V}$ , two vertices  $X, Y \in \mathbf{V}$  are  $d$ -separated by a set of vertices  $\mathbf{Z} \subset \mathbf{V}$ , written  $(X \perp Y | \mathbf{Z})$ , if and only if  $\mathbf{Z}$  blocks all paths between  $X$  and  $Y$ .*

$d$ -separation is a relation among three sets of variables in a DAG. If one is interpreting DAGs as causal networks (as in this article), then  $d$ -separation must represent a *causal* relation among the three sets of variables. By contrast, conditional independence represents a *statistical* relation among them. One might say that  $\mathbf{X}$  is *causally screened off* from  $\mathbf{Y}$  given  $\mathbf{Z}$  whenever  $\mathbf{X}$  is  $d$ -separated from  $\mathbf{Y}$  given  $\mathbf{Z}$ . Of course, the significance of this causal relation is found in the statistical distributions that can be supported by the causal structure. A set of variables  $\mathbf{X}$  is  $d$ -separated from the set  $\mathbf{Y}$  given the set  $\mathbf{Z}$  in a causal structure if and only if for all probability distributions over the causal structure,  $\mathbf{X}$  is conditionally independent of  $\mathbf{Y}$  given  $\mathbf{Z}$ .